

## SNP Resources: Finding SNPs Discovery and Databases

Mark J. Rieder, PhD

NIEHS Variation Workshop

January 30-31, 2006

## SNP Resources: SNP discovery and cataloging

1. SNP discovery/genotyping: Genome-wide approaches
  - ✓ SNP Consortium
  - ✓ HapMap
2. The current state of SNP resources
3. Comprehensive SNP discovery
  - NIEHS SNPs - Environmental Genome Project

SNP Databases - **How to Manual for finding SNPs**  
In class - Tutorial

## Genetic Markers: Overview

1. RFLPs (SNPs circa 1980)
2. Microsatellites (SSLP; di-, tri-, tetranucleotide repeats)
  - 1/50,000 bp
  - Linkage Studies 300-400 markers (~1 Mbp)
  - Multi-allelic/High heterozygosity/informative
  - Complex genotyping assays
3. Single Nucleotide Polymorphisms (SNPs)
  - Most frequent genetic variant (base substitutions)
  - 1/1000 bp (comparing randomly selected chromosomes)
  - Biallelic/less informative
  - Simplified genotyping platforms (+/- calling)

## Development of a genome-wide SNP map: How many SNPs?

Table 1 • Occurrence of SNPs in the human population

Minimal allele frequency	Expected SNP number (millions)	Expected SNP frequency (bp)
1%	11.0	290
5%	7.1	450
10%	5.3	600
20%	3.3	960
30%	2.0	1,570
40%	0.97	3,280

Nickerson and Kruglyak, *Nature Genetics*, 2001

~ 10 million common SNPs (> 1- 5% MAF) - 1/300 bp

How has SNP discovery progressed toward this goal

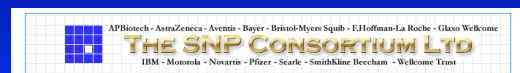
## Finding SNPs: Marker Discovery and Methods

SNP discovery has proceeded in two distinct phases:

- 1 - SNP Identification
  - Define the alleles
  - Map this to a unique place in the genom
- 2 - SNP Characterization
  - Determination of the genotype in many individuals
  - Population frequency of SNPs

## Finding SNPs: Marker Discovery and Methods

SNP Discovery has proceeded in two distinct phases: 1 - **SNP Discovery\*\***/Characterization



2 - SNP Discovery/Characterization\*\*

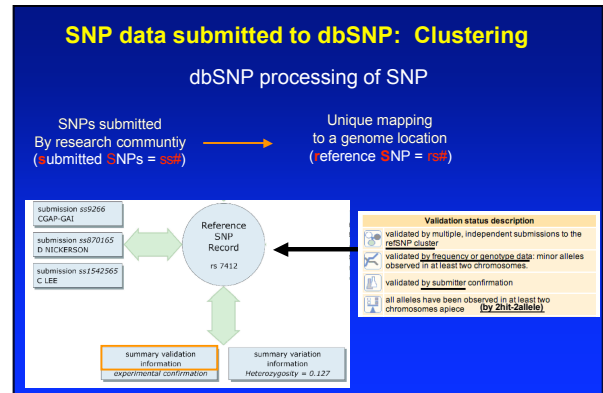




## SNP Discovery: dbSNP database

dbSNP  
-NCBI SNP database

The screenshot shows the NCBI dbSNP homepage. It includes a search bar, navigation links for various databases, and a section for dbSNP search options. The search options include fields for ID, description, and date. There is also a section for submission information, including details about the submission process and the types of data accepted.



## Finding SNPs: Marker Discovery and Methods

SNP Discovery has proceeded in two distinct phases:

- 1 - SNP Identification\*\*/Discovery
- 2 - SNP Discovery/Characterization\*\*

The slide shows the logos of 'THE SNP CONSORTIUM LTD' and the 'International HapMap Project'. The SNP Consortium Ltd logo includes the names of its members: APBioscience, Amgen, Aventis, Bayer, Bristol-Myers Squibb, Eli Lilly, Glaxo Wellcome, IBM, Motorola, Novartis, Pfizer, Sanofi-Schering Plough, and Wellcome Trust. The International HapMap Project logo includes the text 'Home | About the Project | Data | Publications'.

## HapMap Project Proposed: Map more SNPs and genotype

The screenshot shows the International HapMap Project website. It includes a header with the project name and a navigation bar. Below the header, there is a list of 'Participating Groups' from various institutions around the world. The main content area lists the project's goals and objectives.

- Increase SNP density over the first 6 - 12 months
- Ultimately produce a fine scale genetic map (HapMap) which would serve as a common resource for all biomedical researchers
- Genotype 600,000 SNPs genome-wide
- Four populations: CEPH (Europe), Yoruban (Africa), Japanese/Chinese (Asian)

## HapMap SNP Discovery: Prior to Genotyping

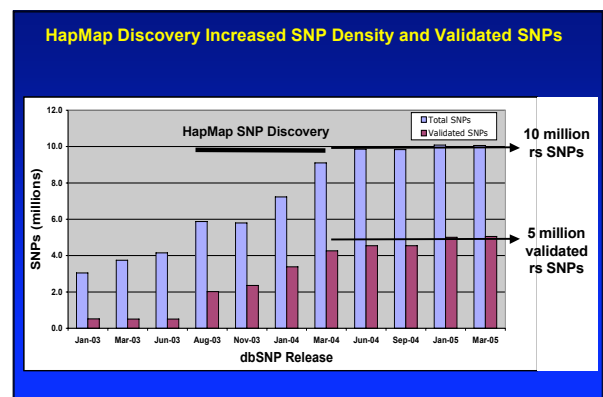
Initiation of project planning (July 2001):  
2.8 million SNPs (1.4 million validated) - 1/1900 bp  
Nov 2003 - 5.7 million (2 million validated) - 1/1500 bp  
Feb 2004 - 7.2 million (3.3 million validated) - 1/900 bp

**Generate more SNPs:**  
Random Shotgun Sequencing

Genomic DNA (multiple individuals) → Sequence and align (reference sequence)

**Other Sources of SNPs:**  
Perlegen (Affymetrix chips) SNP data (chr22)  
Sequence chromatograms from Celera project

TACGCCATA TCAGGAGAT  
GTTACGCCAATACAGGATCCAGGAGATTACC Draft Human Genome



## Development of a genome-wide SNP map: How many SNPs?

Minimal allele frequency	Expected SNP number (millions)	Expected SNP frequency (bp)
1%	11.0	290
5%	7.1	450
10%	5.3	600
20%	3.3	960
30%	2.0	1,570
40%	0.97	3,280

Nickerson and Kruglyak, Nature Genetics, 2001

~ 10 million common SNPs (> 1- 5% MAF) - 1/300 bp

Feb 2001 - 1.42 million (1/1900 bp)

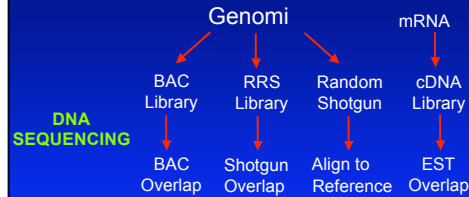
Nov 2003 - 2.0 million (1/1500 bp)

Feb 2004 - 3.3 million (1/900 bp)

Mar 2005 - 5.0 million (validated - 1/600 bp)

When will we have them all?

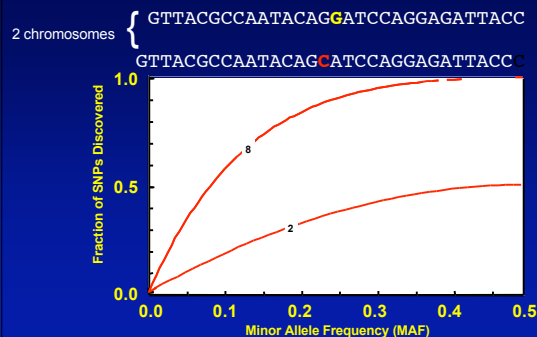
## Finding SNPs: Sequence-based SNP Mining



## RANDOM Sequence Overlap - SNP Discovery

GTTACGCCAATACAGGATCCAGGAGATTACC  
GTTACGCCAATACAGCATCCAGGAGATTACC

## SNP discovery is dependent on your sample population size



## SNP Characterization/Genotyping

Minimal allele frequency	Expected SNP number (millions)	Expected SNP frequency (bp)
1%	11.0	290
5%	7.1	450
10%	5.3	600
20%	3.3	960
30%	2.0	1,570
40%	0.97	3,280

Nickerson and Kruglyak, Nature Genetics, 2001

~ 10 million common SNPs (>1- 5% MAF) - 1/300 bp

Mar 2005 - 5.0 million (validated/mapped - 1/600 bp)

5.0/10.0 = 50% of all common SNPs (validated)!

## HapMap Project Proposed: Map more SNPs and genotype



Participating Groups	Participating Groups
Baylor College of Medicine (USA)	Johns Hopkins School of Medicine (USA)
Beijing Genomics Institute (China)	McGill University & Genome Quebec Innovation Centre (Canada)
Beijing Normal University (China)	ParkeBio Biotechnology (USA)
Broad Institute of Harvard and MIT (USA)	Paragon Science (USA)
Center for Statistical Genetics, University of Michigan (USA)	Riken (Japan)
Chinese National Human Genome Center at Beijing (China)	The Chinese University of Hong Kong (China)
Chinese National Human Genome Center at Shanghai (China)	The University of Hong Kong (China)
Cold Spring Harbor Laboratory (USA)	University of California, San Francisco (USA)
Eubios Ethics Institute (Japan)	University of Basel (Nigeria)
Health Sciences University of Hokkaido (Japan)	University of Oxford (UK)
Hong Kong University of Science and Technology (China)	University of Tokyo (Japan)
Howard University (USA)	University of Utah (USA)
Illumina (USA)	Washington University, St. Louis (USA)
	Wellcome Trust Sanger Institute (UK)

• Genotype 600,000 SNPs genome-wide

• Four populations:

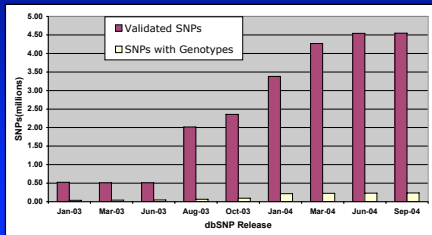
- CEPH (CEU) (Europe - n = 90, trios)
- Yoruban (YRI) (Africa - n = 90, trios)
- Japanese (JPT) (Asian - n = 45)
- Chinese (HCB) (Asian - n = 45)

## Finding SNPs: Genotype Data Adds Value to SNPs HapMap Genotyping

- ✓ Confirms SNP as 'real' and 'informative'
- ✓ Minor Allele Frequency (MAF) - common or rare
- ✓ MAF in different populations
- ✓ Detection of SNP x SNP correlations (Linkage Disequilibrium)
- ✓ Determine haplotypes



## Few SNPs in dbSNPs had Genotype Data



## Perlegen Large-scale Genotyping Capacity

### Whole-Genome Patterns of Common DNA Variation in Three Human Populations

David A. Hinds,<sup>1</sup> Laura L. Stuve,<sup>1</sup> Geoffrey B. Nilsen,<sup>1</sup>  
Eran Halperin,<sup>2</sup> Eleazar Eskin,<sup>3</sup> Dennis G. Ballinger,<sup>1</sup>  
Kelly A. Frazer,<sup>1</sup> David R. Cox<sup>1\*</sup>

18 FEBRUARY 2005 VOL 307 SCIENCE

**1.58 millions SNPs genotyped**  
**71 individuals from 3 American populations**  
**European, African and Asian ancestry**

## HapMap Completion

### A haplotype map of the human genome

The International HapMap Consortium\*

Inherited genetic variation has a critical but as yet largely uncharacterized role in human disease. Here we report a public database of common variation in the human genome: more than one million single nucleotide polymorphisms (SNPs) for which accurate and complete genotypes have been obtained in 269 DNA samples from four populations, including ten 500-kilobase regions in which essentially all information about common DNA variation has been extracted. These data document the generality of recombination hotspots, a block-like structure of linkage disequilibrium and low haplotype diversity, leading to substantial correlations of SNPs with many of their neighbours. We show how the HapMap resource can guide the design and analysis of genetic association studies, shed light on structural variation and recombination, and identify loci that may have been subject to natural selection during human evolution.

**Nature - Oct 27 (2005)**

- 2005-06-01: HapMap public release #1b.1  
This is the final Phase I data freeze as used in analyses for the upcoming primary HapMap publication (see [Data freezes](#) for more info). Also, note that with this release the abbreviation for the Han Chinese in Beijing population is changed to CHB. (See [Guidelines for Referring to HapMap Populations](#) for more info.)  
Summary of genotyped SNPs:

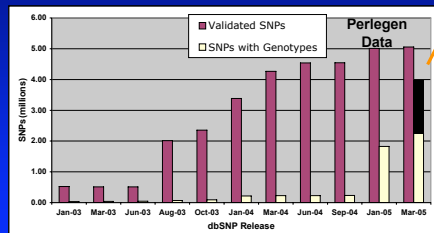
Populations	CEU	CHB	JPT	YRI
Genotyped SNPs	1,105,072	1,098,689	1,099,420	1,078,451

- 2005-10-24: HapMap Public Release #19  
Genotypes, frequencies and assays for phase I and phase II of the HapMap project are now available for [bulk download](#). The files contain all phase I and II data combined.

Populations	CEU	CHB	JPT	YRI
Total QC SNPs	3,957,458	3,953,524	3,952,623	3,855,955
Total Genotyped SNPs	3,954,684	3,912,290	3,912,290	3,857,458

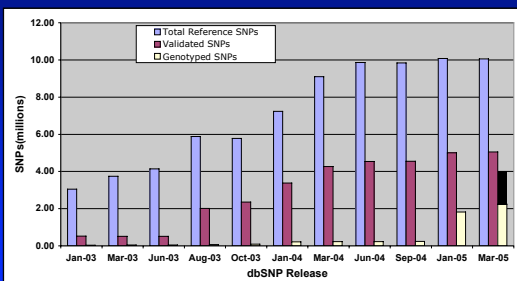
**HapMap + Perlegen**

## dbSNP: Increasing numbers of SNPs now have genotype data



HapMap  
Phase II  
Perlegen

## Current State of dbSNP



**Many SNPs left to validate and characterize.**

## Increasing SNP Density: HapMap ENCODE Project

ENCODE ENCYclopedia Of DNA Elements

Catalog all functional elements in 1% of the genome (30 Mb)

10 Regions x 500 kb/region (Pilot Project)

David Altshuler (Broad), Richard Gibbs (Baylor)

16 CEU, 16 YRI, 8 HCB, 8 JPT

Comprehensive PCR based resequencing across these regions

ENCODE Region Information

Project

Resequencing

Project

Genotyping

Perlegen

Genotyping

Component

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE Link

Link

ENCODE

## Development of a genome-wide SNP map: How many SNPs?

**Table 1 • Occurrence of SNPs in the human population**

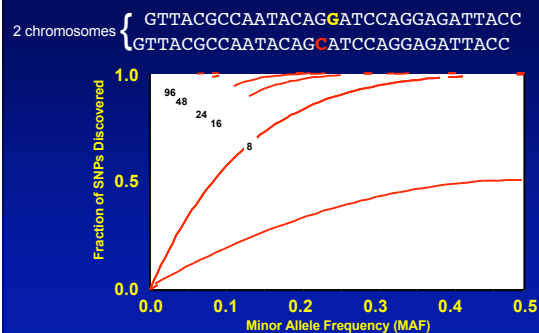
Minimal allele frequency	Expected SNP number (millions)	Expected SNP frequency (bp)
1%	11.0	290
5%	7.1	450
10%	5.3	600
20%	3.3	960
30%	2.0	1,570
40%	0.97	3,280

Nickerson and Kruglyak, Nature Genetics, 2001

~ 10 million common SNPs (>1- 5% MAF) - 1/300 bp  
 Mar 2005 - 5.0 million (validated - 1/600 bp)

~4.0 million validated SNPs with genotypes!  
 (HapMap confirmed, allele frequency/population, SNPxSNP correlations (LD), haplotypes)

## SNP discovery is dependent on your sample population size



National Institute of Environmental Health Sciences  
 Environmental Genome Project  
**NIEHS SNPs**

Search Site  Go

**Goal:** Comprehensively identify all common sequence variation in candidate genes

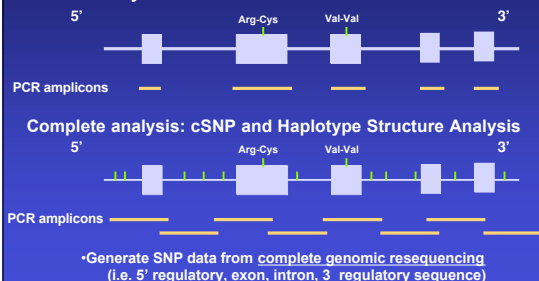
**Initial biological focus:** Candidate environmental response genes involved in DNA repair, cell cycle, apoptosis, metabolism, cell signaling, and oxidative stress.

**Approach:** Direct resequencing of genes

**Samples:** PDR 90 ethnically diverse individuals representative of U.S. population (397 genes)  
 EGP95 95 samples from 4 ethnic groups (23 HapMap Asians, 22 HapMap Europeans, 15 HapMap Yorubans, 12 African Americans, 24 Hispanic ) (170 genes)

## Targeted SNP Discovery

### Directed analysis: cSNPs



## Summary of NIEHS SNP genotypes in dbSNP

**Table 1. Summary of genotype data contained in dbSNP**

Data set	Genotypes	SNPs	Populations	Individuals	Average SNP density	Reference
HAPMAP	159,842,776	954,302	4	270	3149	(International HapMap Consortium 2003)
PERLEGEN	110,385,051	1,576,578	3	71	1938	(Hindis et al. 2003)
Allymetrix	6,189,466	125,778	6	116	24,029	(Kennedy et al. 2003)
TYC	4,932,382	19,048	12	1963	312,254	(International SNP Map Working Group 2001)
PCA	2,164,102	12,635	12	99	77,445	(Crawford et al. 2004)
PCA/UW	275,194	15,981	2	47	153,861	(Crawford et al. 2004)
IPCA	176,162	3801	3	47	430,361	(Immune Immunity PCA, http://innateimmunity.net/)
NIH/PDR	159,549	1982	1*	448	1,419,125	(Collins et al. 1998)
WICAR	33,240	1462	1	130	2,011,277	(Freudenberg-Hua et al. 2003)
HG_BONN	24,522	320	1	143	5,284,550	(Freudenberg-Hua et al. 2003)

\*The NIH/PDR data contains a single mixed population.

Current numbers  
 554 genes sequenced  
 12.76 Mb scanned  
 75,580 genotyped SNPs identified  
 7 million genotypes deposited in dbSNP

Nov 2005 -Zaitlen et al. Genome Research 15:1594-1600

## Development of a genome-wide SNP map: How many SNPs?

**Table 1 • Occurrence of SNPs in the human population**

Minimal allele frequency	Expected SNP number (millions)	Expected SNP frequency (bp)
1%	11.0	290
5%	7.1	450
10%	5.3	600
20%	3.3	960
30%	2.0	1,570
40%	0.97	3,280

Nickerson and Kruglyak, Nature Genetics, 2001

~ 10 million common SNPs (>1- 5% MAF) - 1/300 bp

NIEHS SNPs = 1/180 bp (n = 95, 4 pops)

HapMap ENCODE = 1/160 (n = 48, 3 pops)

Comprehensive resequencing can identify the vast majority of SNPs in a region

## SNP Discovery: dbSNP database

dbSNP (Perlegen/HapMap)

Minor Allele Freq. (MAF)

NIEHS SNPs

SNP Distribution

Number

Minor Allele Frequency

60%

25%

15%

Rarer and population specific SNPs are found by resequencing

**PDR** = 90 ethnically diverse individuals representative of U.S. population (397 genes - ~55,000 SNPs )

## HapMap Populations

## Non-HapMap Populations

- | Array (1536) | Site Conversion Rate (%) | Average Site Coverage (%) | Concordance (%) |
|--------------|--------------------------|---------------------------|-----------------|
| 1            | 85                       | 96.6                      | 99.7            |
| 2            | 91                       | 97.7                      | 99.5            |
| 3            | 82                       | 98.5                      | 99.3            |

[illegible]

- SNPs have been rapidly adopted as the genetic marker of choice.
- Approximately 10 million common SNPs exist in the human genome (1/300 bp).
- Random SNP discovery processes generate many SNPs (TSC and HapMap).
- Random approaches to SNPs discovery have reached limits of discovery and validation (1/600 bp; 50% SNP validation)
- Most validated SNPs (5 million) will be genotyped by the HapMap (3 pops)
- Resequencing approaches continue to catalog important variants (rarer)
- NIEHS SNPs has generated SNP data on >550 candidate genes and 75 K SNPs